

E³ : Keyphrase based News Event Exploration Engine

Nikita Jain Swati Gupta Dhaval Patel
 Department of Computer Science and Engineering
 Indian Institute of Technology Roorkee, India
 nk27jain@gmail.com, {sg123.dcs2014, patelfec}@iitr.ac.in

ABSTRACT

This paper presents a novel system E³ for extracting keyphrases from news content for the purpose of offering the news audience a broad overview of news events, with especially high content volume. Given an input query, E³ extracts keyphrases and enrich them by tagging, ranking and finding role for frequently associated keyphrases. Also, E³ finds the novelty and activeness of keyphrases using news publication date, to identify the most interesting and informative keyphrases.

1. INTRODUCTION

News media are publishing ideas, events and opinions in an increasingly wide range of data formats such as news articles, headlines, videos, tweets, hashtags and others. The explosion of Big news data has sparked the text and data mining research communities to focus on developing systems for news data exploration and analysis. Broadly, two types of news data exploration systems are developed till date: *Event centric* (GDELT [1], EventRegistry [2]) and *Content centric* (STICS [3], EMM [4]). In Event centric system, input query maps to real world events, whereas, the content centric system outputs related news articles of a given query.

Although, both types of systems provide up-to-date news information in real time, but they overload the user with the large amounts of results. For instance, given input query “2014 FIFA World Cup”, event centric EventRegistry suggested 11,504 news events, and the content centric STICS suggested 1,286,369 news articles having multiple organizations, people and places mentioned. Clearly, there is a need of a system that enables readers to get a broad overview of the news data generated in response of user query.

In this paper, we propose a keyphrase based news exploration engine E³ to summarize high volume news data. In our context, keyphrase is a short and meaningful chunk of text that describes an important news concepts, news entities, etc. For instance, “Bihar election”, “Bihar bjp” are examples of keyphrases. Our proposed work is keyphrase centric as recent literature has shown that Keyphrase min-

ing is able to generate a wide range of informational and important phrases from large documents (KEA [5], Micro-gram [6], ToPMine [7], and SegPhrase [8]). As shown in Figure 1, engine E³ works in two phases: keyphrase extraction followed by keyphrase enrichment. Keyphrase extraction is performed on multi-form news (like article, keywords, and others). As existing keyphrase mining approaches performs poorly on news data, we propose a novel keyphrase extraction technique that leverages linguistic-syntactic feature, and performs very well on short texts like news headlines and video captions. The keyphrase enrichment phase finds important and interesting information related to the keyphrase, such as connected entities, novel (emerging), and active news concepts, and the role played by most frequent entities present in set of keyphrases.

In summary, given an input query q , engine E³ generates a keyphrase template, as shown in Figure 2, where keyphrases are organized into three sections: Type Discovery, Keyphrase Ranking and InfoBox Miner. The keyphrase of type Person, Location and Organization are kept in Type section. The Novel, Active and Frequent keyphrases are annotated in Keyphrase Ranking section. Information like Role played, top most associated Types and Keyphrase for a selected keyphrase is stored in InfoBox section.

2. SYSTEM OVERVIEW

Figure 1 gives an overview of E³ architecture.

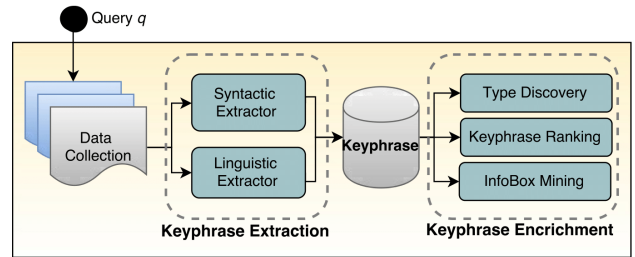


Figure 1: E³ System Architecture

2.1 Data Collection

For input query q , our data collection module prepares related news data published by news media. The module can either use news search engines like Google News, Yahoo News or online news structured repositories like GDELT [1], EMM [4], IMM [9]. We use our in-house IMM system which periodically extracts news headlines (video title) along with their URL, publication date and meta-keywords. To prepare

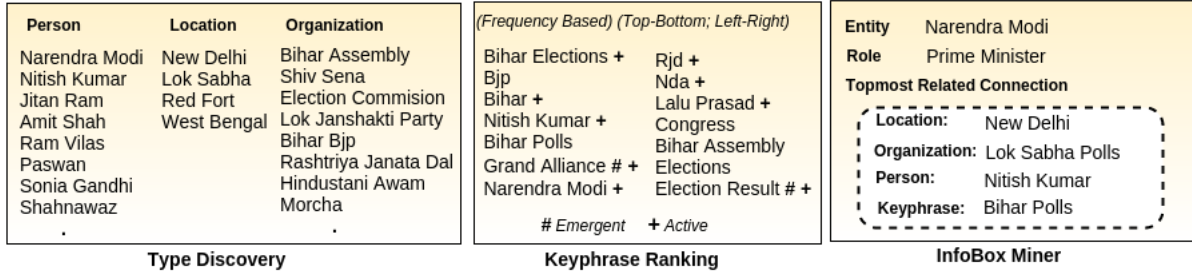


Figure 2: E^3 System Working Example for q : “Bihar Election”

news data related to q , we select URL, if URL’s headlines (video title) or URL’s meta-keywords contains q .

In summary, for given query q , we prepare a dataset R containing several news records R_q , where each record is described by *Quintuple* {Headline, Keywords, Meta-description, Article, Publication Date}. For “Bihar Election” query, we retrieved 216 records from iMM for further processing.

2.2 Keyphrase Extraction

Next, we extract keyphrases from all the records of news data R . As meta-keywords are small group of meaningful words, a naive solution is to output the meta-keywords as keyphrases. However, not all records in R contain meta-keyword. As a result, meta-keywords are not sufficient enough to describe the news data completely. For example, around 20% news headlines, obtained for “Bihar Election”, do not have a meta-keyword. Hence, we require an efficient selection of the relevant phrases.

We observed that news headlines are short in length and contains special tokens such as colon, apostrophe, quotes, hash, dash to emphasize important information. On the other hand, meta-descriptions and news articles are long passage texts and are governed by grammatical rules. Thus, we propose two different keyphrase extractors to handle both kinds of writing styles.

- **Syntactic Extractor** utilises special characters such as colon (:), apostrophe (’), quotes (“, ’), hash (#), dash (-) for keyphrase extraction. In case of colon (dash), the news headline is tokenized into two parts using colon (dash) and both parts are declared as keyphrases. In case of quotes (hash) the part of text enclosed inside the quotes and hashtag containing the hash is declared as keyphrase respectively.
- **Linguistic Extractor** applies language specific part of speech (POS) tagging on meta-description and article and then annotate collocated nouns, adjectives, noun apostrophe (’) connector and numbers present in the input text. These annotated tokens are further used as keyphrases.

At the end, when all the news records in R are processed, we obtain a set of keyphrases R_k , along with the number of times they are generated. For “Bihar Election” query, we obtained around 6000 keyphrases for further processing.

2.3 Keyphrase Enrichment

The size of generated keyphrases R_k may be large and noisy. To resolve this problem, our Keyphrase Enrichment module helps in extracting valuable and actionable information by filtering and ranking the extracted keyphrases. The keyphrases are filtered using news media specific stop-words such as update, video, photo, pti and others. Next, we

apply case normalization and remove duplicate keyphrases. At this point noisy keyphrases are removed. The remaining keyphrases are passed through the Type discovery, Keyphrase ranking and InfoBox mining modules.

- In **Type Discovery** module, NER tagger is used to classify keyphrase into three types: *Person*, *Location* and *Organization*. As existing NER taggers do not perform well on Indian named entities, we use a separate list¹ for Indian named entities, prepared through in-house research work. A keyphrase without any above NER type, are termed as a *News Concept*. For “Bihar Election” query, a sample keyphrases for each type is shown in Figure 2.
- Using **Keyphrase Ranking** module, keyphrases are organized according to the value of frequency, novelty and activeness. The frequency of the keyphrase is already computed during keyphrase extraction process. To compute the value of novelty and activeness, we first extracts the time intervals q_t of an input query q , during which the q was highly popular in news headlines. Next, a keyphrase is *novel* (denoted by # in Figure 2) if its frequency is very high in news headlines only during q_t . Similarly, a keyphrase is *active* (denoted by + in Figure 2) if its frequency is very high around q_t . For instance, “Grand Alliance” and “NDA” are discovered as novel and active keyphrases respectively for “Bihar Election” query.
- **InfoBox Miner** discovers personalized information for selected keyphrase. The InfoBox displays role of keyphrase k that it played with respect to query q . The type-wise top most connections k have, determined with help of co-occurrence value of the keyphrases in R_k . A phrase frequently located near k in the collected news corpus is labeled as the k ’s role. Generally, keyphrases with type person and organization are preferred for InfoBox mining. Figure 2 shows InfoBox for entity “Narendra Modi” for query “Bihar Election”.

3. SUMMARY

The engine is tested for varying the input query ranging from general topics (e.g., Election, ISIS) to specific topics (e.g., Paris Attack, Gravitational Wave) and compared the results² with KEA, ToPMine and Micro-ngram. We found that our system outperforms existing approaches in terms of quality and quantity of keyphrases generated. As our engine is online, we can demonstrate the working to the conference participants.

¹<https://github.com/NikkiJain09/Transliteration>

²goo.gl/yoLXTh

4. REFERENCES

- [1] K. Leetaru and P. A. Schrodt, *GDELT: Global data on events, location, and tone*. International Studies Association Annual Convention, 2013.
- [2] G. Leban, B. Fortuna, J. Brank, and M. Grobelnik, *Event Registry: Learning About World Events from News*. World Wide Web, 2014.
- [3] J. Hoffart, D. Milchevski, and G. Weikum, *STICS: Searching with Strings, Things, and Cats*. Special Interest Group on Information Retrieval, 2014.
- [4] R. Steinberger, B. Pouliquen, and E. V. der Goot, *An introduction to the Europe Media Monitor family of applications*. Special Interest Group on Information Retrieval, 2009.
- [5] O. Medelyan and I. H. Witten, *Thesaurus Based Automatic Keyphrase Indexing*. Joint Conference on Digital Libraries, 2006.
- [6] K. Wang, C. Thrasher, E. Viegas, X. Li, and B.-j. P. Hsu, *An Overview of Microsoft Web N-gram Corpus and Applications*. North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, 2010.
- [7] A. El-Kishky, Y. Song, C. Wang, C. R. Voss, and J. Han, *Scalable Topical Phrase Mining from Text Corpora*. Very Large Data Bases, 2014.
- [8] J. Liu, J. Shang, C. Wang, X. Ren, and J. Han, *Mining Quality Phrases from Massive Text Corpora*. Special Interest Group on Management of Data, 2015.
- [9] S. Mazumder, B. Bishnoi, and D. Patel, *News Headlines: What They Can Tell Us?* IBM Collaborative Academia Research Exchange, 2014.